

TWEETEVAL: Unified Benchmark and Comparative Evaluation for Tweet Classification

Francesco Barbieri[♣] Jose Camacho-Collados[†]
Leonardo Neves[♣] Luis Espinosa-Anke[†]

[♣]Snap Inc., Santa Monica, CA 90405, USA

[†]School of Computer Science and Informatics, Cardiff University, United Kingdom

[♣]{fbarbieri, lneves}@snap.com,

[†]{camachocolladosj, espinosa-ankel}@cardiff.ac.uk

Abstract

The experimental landscape in natural language processing for social media is too fragmented. Each year, new shared tasks and datasets are proposed, ranging from classics like sentiment analysis to irony detection or emoji prediction. Therefore, it is unclear what the current state of the art is, as there is no standardized evaluation protocol, neither a strong set of baselines trained on such domain-specific data. In this paper, we propose a new evaluation framework (TWEETEVAL) consisting of seven heterogeneous Twitter-specific classification tasks. We also provide a strong set of baselines as starting point, and compare different language modeling pre-training strategies. Our initial experiments show the effectiveness of starting off with existing pre-trained generic language models, and continue training them on Twitter corpora.

1 Introduction

Modern NLP systems are typically ill-equipped when applied to noisy user-generated text. The high-paced, conversational and idiosyncratic nature of social media, paired with platform-specific restrictions (e.g., Twitter’s character limit), requires tackling additional challenges, for example, POS tagging (Derczynski et al., 2013), lexical normalization (Han and Baldwin, 2011; Baldwin et al., 2015), or named entity recognition (Ritter et al., 2011; Baldwin et al., 2013). In other more generic contexts, these challenges can be considered solved or are simply non-existent. Moreover, other apparently simple tasks such as sentiment analysis have proven to be hard on Twitter data (Poria et al., 2020), among others, due to limited amount of contextual cues available in short texts (Kim et al., 2014). In addition to these and other inherent difficulties, advances in NLP for user-generated data are hindered by its highly fragmented landscape

and the lack of a unified evaluation framework. In the current era of pretraining and Language Models (LMs), this is particularly relevant, as these models exhibit a versatility that currently cannot be gauged comparably across Twitter datasets and tasks. This is not the case, however, in more ordinary textual genres and domains. For instance, well known benchmarks like SentEval (Conneau and Kiela, 2018), GLUE (Wang et al., 2019b) and SuperGLUE (Wang et al., 2019a) include standard NLP tasks such as language inference, paraphrase detection or sentiment analysis, among others. It is undisputable that these benchmarks have contributed to the fast development of language understanding techniques, and LMs in particular, as they have enabled comprehensive evaluations across several tasks in fair and reproducible conditions.

We thus take inspiration from the above to develop TWEETEVAL, a benchmark for tweet classification in English. TWEETEVAL is a standardized test bed for seven tweet classification tasks. These are: sentiment analysis, emotion recognition, offensive language detection, hate speech detection, stance prediction, emoji prediction, and irony detection. We develop a unified framework, unified criteria for train/validation/test splits, and evaluate strong baselines inspired by current SotA in these tasks. We also evaluate transformer-based models, trained entirely and partially on Twitter data, with which we aim to establish a competitive high bar for subsequent contributions. The contributions of this paper are therefore as follows: (1) we compile, curate and release a suite of tasks under the umbrella of a new benchmark: **TWEETEVAL¹**, a unified framework comprising several tweet classification tasks; and (2) we **evaluate state-of-the-art LMs in this new framework**, and shed light on the effect of training with different corpora.

¹The unified TWEETEVAL benchmark is available at: <https://github.com/cardiffnlp/tweeteval>


Dataset	Tweet	Label
Emoji	Thx for showing this newbie passholder around @ Disneyland	
Emotion	I love swimming for the same reason I love meditating...the feeling of weightlessness.	joy
Hate	Another illegal alien that shouldn't be in America killed an innocent American couple! #BuildThatWall	hateful
Irony	Leaving whilst its dark is fun. #not	ironic
Offensive	Are we all ready to sit and watch Indakurate Passcott play football?	non-offensive
Sentiment	Hmmmm where are the #BlackLivesMatter when matters like this a rise... kids are a disgrace!!	negative
Stance(<i>fem</i>)	Rather be an "ugly" feminist then be these sad people that throws hat on people that believes in equality!	in favour

Table 1: Tweet samples for each of the tasks we consider in TweetEval, alongside their label in their original datasets. We use (*fem*) to refer to the *feminism* subset of the stance detection dataset.

2 TweetEval: The Benchmark

In this section, we describe the compilation, curation and unification procedure behind the construction of TWEETVAL and its corresponding tasks, as well as relevant statistics and evaluation metrics. We also show, in Table 1, a sample tweet and its corresponding label from the original task.

2.1 Tasks

Emotion Recognition. This task consists of recognizing the emotion evoked by a tweet. We use the dataset of the most participated task of SemEval2018, "Affects in Tweets" (Mohammad et al., 2018). The original competition was framed as a multi-label classification problem, including 11 emotions. The integration into TWEETVAL consists of re-purposing this multi-label dataset into multi-class classification, keeping only the tweets labeled with a single emotion. Since the amount of tweets with single labels was scarce, we selected the most common four emotions (Anger, Joy, Sadness, Optimism)².

Emoji Prediction. This task consists in, given a tweet, predicting its most likely emoji, and is based on the Emoji Prediction challenge at Semeval2018 (Barbieri et al., 2018). It only considers tweets with one emoji (irrespective of its position), which is used as classification label. The test set is the same as in the original publication, but we limit the training and validation splits to 50,000 tweets, in order to comply with Twitter distribution policies. The label set comprises 20 different emoji, and due to their skewed distribution, this task proved to be highly difficult, with low overall numbers. Specifically, more than 42% of the tweets are labeled with the 3 most frequent emoji (❤️, 😊, and 🤔).

²We selected those emotions with a minimum frequency of 300 examples in the training set.

Task	Lab	Train	Val	Test
Emoji prediction	20	45,000	5,000	50,000
Emotion det.	4	3257	374	1421
Hate speech det.	2	9,000	1,000	2,970
Irony detection	2	2,862	955	784
Offensive lg. id.	2	11,916	1,324	860
Sent. analysis	3	45,389	2,000	11,906
Stance detection	3	2620	294	1249
Stance/Abortion	3	587	66	280
Stance/Atheism	3	461	52	220
Stance/Climate	3	355	40	169
Stance/Feminism	3	597	67	285
Stance/H. Clinton	3	620	69	295

Table 2: Number of labels and instances in training, validation, and test sets for each dataset. The specific statistics of each target domain in the stance detection task is included at the bottom.

Irony Detection. This task consists of recognizing whether a tweet includes ironic intents or not. We use the Subtask A dataset of the SemEval2018 Irony Detection challenge (Van Hee et al., 2018). Note that this dataset was artificially balanced to make the task more accessible.

Hate Speech Detection. This task consists in predicting whether a tweet is hateful or not against any of two target communities: immigrants and women. Our dataset of choice stems from the SemEval2019 Hateval challenge (Basile et al., 2019).

Offensive Language Identification. This task consists in identifying whether some form of offensive language is present in a tweet. For our benchmark we rely on the SemEval2019 Offenseval dataset (Zampieri et al., 2019).

Sentiment Analysis. The goal for the sentiment analysis task is to recognize if a tweet is positive, negative or neutral. We use the Semeval2017 dataset for Subtask A (Rosenthal et al., 2019),

which includes data from previous runs (2013, 2014, 2015, and 2016) of the same SemEval task.

Stance Detection. Stance detection is the task to determine, given a piece of text, whether the author has a favourable, neutral, or negative position towards a proposition or target. We use the SemEval2016 shared task on Detecting Stance in Tweets (Mohammad et al., 2016). In the original task, five target domains are given: abortion, atheism, climate change, feminism and Hillary Clinton. Unlike the other tasks, training is provided separately for each target domain, which we use to extract individual validation sets.

2.2 Statistics and evaluation metrics

Table 2 includes the TWEETVAL datasets statistics after unification.³ Data sizes range from a few hundred instances for training to over 40,000. Note that the preprocessing pipeline is equal for all tasks: user mentions are anonymized and line breaks and website links are removed.

Evaluation metrics. We use the same evaluation metric from the original tasks, which is macro-averaged F1 over all classes, in most cases. There are three exceptions: stance (macro-averaged of F1 of favor and against classes), irony (F1 of ironic class), and sentiment analysis (macro-averaged recall). Similar to GLUE (Wang et al., 2019b), we also introduce a global metric (TE) based on the average of all dataset-specific metrics.

3 Language Models for Tweet Classification

Transformer-based LMs such as GPT (Radford et al., 2018), BERT (Devlin et al., 2019) or XLNET (Yang et al., 2019) have taken the NLP field by storm, outperforming previous linear models and neural network methods based on LSTMs or CNNs in many tasks, including sentence and text classification (Wang et al., 2019b).

The functioning of these language models for tweet classification is conceptually simple. First, they are trained on a large unlabeled corpus. Then, they are fine-tuned to the task for where an appropriate training set exists. For social media text, however, one may question whether existing pre-trained models trained on standard corpora are optimal. We thus compare three different strategies

³The validation sets are randomly sampled from the training set for those tasks where no validation split is provided in the original dataset.

which differ in the training data: (1) Using an existing large pre-trained LM; (2) using an existing architecture, but training from scratch using only Twitter data; and (3) starting with an original pre-trained LM and continue to train with Twitter data, keeping the original tokenizer and the same masked LM loss.

We consider these three techniques as we are interested in exploring whether a Twitter-specific LM should be trained on Twitter only or if it should be initialized with weights learned during pre-training on standard corpora, *and then* be trained on Twitter. The latter option has indeed three theoretical advantages: (1) these models are generally trained on large amounts of text corpora, and reproducing the same experiment would be extremely expensive even if we had same amount of Twitter data; (2) learning on different types of text corpora make the models more robust and knowledgeable about the world; and (3) some models such as RoBERTa (Liu et al., 2019) or GPT-2 (Radford et al., 2019) are not unfamiliar with internet language and slang, as part of their underlying training corpora contains Reddit data (38GB).

4 Evaluation

4.1 Experimental setting

Neural language model. Among all the available language models we selected RoBERTa (Liu et al., 2019) as it is one of the top performing systems in GLUE. Moreover, it does not employ the Next Sentence Prediction (NSP) loss (Devlin et al., 2018), making the model more suitable for Twitter where most tweets are composed of a single sentence.

Language model pre-training. We use three different RoBERTa variants: pre-trained RoBERTa-base⁴ (RoB-Bs), the same model but re-trained on Twitter (RoB-RT) and trained on Twitter from scratch (RoB-Tw). RoB-RT and RoB-Tw are trained with early stopping on the validation split and learning rate $1.0e^{-5}$. Both models converged after about 8/9 days on 8 NVIDIA V100 GPUs.⁵

Twitter corpus. We train RoB-RT and RoB-Tw on 60M tweets⁶ obtained by extracting a large corpus of English tweets⁷ (using the automatic labeling provided by Twitter). We only considered tweets

⁴RoBERTa-base was trained on 160G of uncompressed text.

⁵We used the Huggingface *transformers* library. The estimated cost for each language model is USD 4,000 on Google Cloud.

⁶584 million tokens (3.6G of uncompressed text).

⁷Crawled with the stream API from May’18 to August’19.

		Emoji	Emotion	Hate	Irony	Offensive	Sentiment	Stance	ALL
Val	SVM	25.0	63.8	73.1	63.4	72.7	68.4	67.9	62.0
	FastText	23.2	62.9	71.7	62.7	70.0	62.2	67.3	60.0
	BLSTM	19.4	62.6	72.1	60.6	72.1	61.9	63.4	58.9
	RoB-Bs	24.7±0.3 (24.3)	73.1±1.7 (74.9)	76.5±0.3 (76.6)	73.7±0.6 (73.7)	77.1±0.6 (77.6)	71.4±1.9 (72.7)	71.4±1.9 (73.9)	67.7
	RoB-RT	24.4±1.5 (26.2)	75.4±1.5 (77.0)	77.8±1.1 (79.6)	74.7±1.5 (75.6)	77.2±0.6 (77.7)	73.0±1.2 (74.2)	72.9±1.0 (75.2)	69.4
	RoB-Tw	23.4±1.1 (24.6)	67.6±0.9 (68.6)	74.3±2.0 (76.6)	70.0±0.3 (70.7)	76.1±0.6 (76.2)	70.5±1.0 (69.4)	68.3±2.4 (71.4)	65.4
Test	SVM	29.3	64.7	36.7	61.7	52.3	62.9	67.3	53.5
	FastText	25.8	65.2	50.6	63.1	73.4	62.9	65.4	58.1
	BLSTM	24.7	66.0	52.6	62.8	71.7	58.3	59.4	56.5
	RoB-Bs	30.9±0.2 (30.8)	76.1±0.5 (76.6)	46.6±2.5 (44.9)	59.7±5.0 (55.2)	79.5±0.7 (78.7)	71.3±1.1 (72.0)	68±0.8 (70.9)	61.3
	RoB-RT	31.4±0.4 (31.6)	78.5±1.2 (79.8)	52.3±0.2 (55.5)	61.7±0.6 (62.5)	80.5±1.4 (81.6)	72.6±0.4 (72.9)	69.3±1.1 (72.6)	65.2
	RoB-Tw	29.3±0.4 (29.5)	72.0±0.9 (71.7)	46.9±2.9 (45.1)	65.4±3.1 (65.1)	77.1±1.3 (78.6)	69.1±1.2 (69.3)	66.7±1.0 (67.9)	61.0
	<i>Best</i>	36.0*	-	65.1	70.5	82.9	68.5	71.0	-
Metric		M-F1	M-F1	M-F1	F ⁽ⁱ⁾	M-F1	M-Rec	AVG (F ^(a) , F ^(f))	TE

Table 3: TweetEval test results. For neural models we report both the average result from three runs and its standard deviation, and the maximum result (parentheses). *Best* results correspond to the best systems in the original shared tasks - they are included for completeness as they not directly comparable. Splits might differ, and * indicates that a larger training set is used. Validation set results are available in the supplemental material.

with at least three tokens and without URLs, as to avoid bot tweets and spam advertising.

Classification fine-tuning. We use the same classification fine-tuning method used in Liu et al. (2019): we add one dense layer to reduce the dimensions of the RoBERTa’s last layer to the number of labels in the classification task, and fine-tune the model on each classification task, training all the parameters simultaneously. We run a minimum parameter search on the starting learning rate ($1.0e^{-3}$, $1.0e^{-4}$, $1.0e^{-5}$, and $1.0e^{-6}$), use early stopping (5 epochs) on the validation set and run each experiment three times with different seeds (1,2,3). Then, we select the highest performing learning rate on the validation set, and use the corresponding model to evaluate on the test set.

Baselines. FastText (Joulin et al., 2017) provides an efficient baseline based on standard features and subword units. We also include an SVM-based baseline with both word and character n-gram features, a model and feature set that has seen great success in recent Twitter-based shared tasks such as emoji prediction (Çöltekin and Rama, 2018) and stance prediction (Mohammad et al., 2018). We finally report the results of a bi-directional LSTM.⁸ Both FastText and the LSTM use 100-dimensional FastText word embeddings (Bojanowski et al., 2017) trained on the 60M Twitter corpus for the

⁸The LSTM has 128 cells, an embedding layer of 100 dimensions, dropout (0.5) and, similarly to the language models, the four learning rate values are tuned in the validation set.

lookup table initialization.

4.2 Results

Table 3 shows the results of all comparison systems on TWEETEval. Perhaps surprisingly, RoBERTa-Base (RoB-Bs) performs well on all tasks, even outperforming the model trained on Twitter data only (RoB-Tw) in most tasks. This can also be attributed to the fact that Twitter is not only noisy text, and formal text can be also found regularly (Hu et al., 2013; Xu, 2017). Using more Twitter data for training might further improve the results of RoB-Tw, but this would also translate into an even more expensive training. However, RoBERTa-Base coupled with additional training on the same Twitter corpus (i.e. RoB-RT) proves more effective.

The only task where a model trained from scratch on Twitter performs better is Irony detection, where RoB-Tw shows to better generalize (RoB-RT F1 drops 13 points from validation to test set, while RoB-Tw F1 5 points). This can be due to two factors: (1) irony used on social media might differ from irony on standard text, (2) tweets in our training data are generally short (79.3 characters on average compared to over 100 characters for most other tasks), and therefore tokenizing the text in less word pieces, and potentially less OOVs, becomes more important to generalize. We note that the low results in the task of emoji prediction (when compared to those obtained in the official SemEval task) are due to the downscaling of the training

data. Because of Twitter’s data distribution policy, at TWEETEVAL we release at most 50k tweets per task, whereas in the original competition, by id sharing, the training data was one order of magnitude bigger. As for the results in the hate speech task, the difference in performance between validation and test set is mainly due to these splits being collected at different timespans, as pointed out by the organizers of the task (Basile et al., 2019). This causes a disparity in topic distribution and thus low performance of the systems optimized towards the validation set.

4.3 Tokenizer analysis

Table 4 includes number of tokens⁹ per tweet for each of the tasks and the difference between word pieces of the pre-trained RoBERTa-base and RoBERTa trained on Twitter from scratch. This comparison is useful to understand if a model recognizes more or less tokens: if the difference between the two RoBERTa tokenizers is high, it means that one model had to split more times a word. We can note that the biggest difference in wordpieces between RoB-Bs and Rob-Tw is 6.8% in the hate detection task. This is expected as these tweets include less standard words, such as insults. On the other hand, except for perhaps emotion detection and offensive language identification, the difference is not significant, considering that the original RoBERTa tokenizer was not trained on Twitter text. Moreover, even if the tokenizer of Rob-RT was not retrained from scratch, this does not mean that Rob-RT could not learn new tokens as they could be learned as sequence of characters during the language modeling re-training phase. This is also the case of emoji, which were not learned in the original RoBERTa model, but BTE includes all their Unicode bytes.

5 Conclusion

We have presented TWEETEVAL, a unified benchmark for tweet classification consisting of seven heterogeneous tasks that are core to social media NLP research. Along with the benchmark, we have included strong baselines as reference, and ran an analysis of LMs with different training strategies. Our results suggest that using a pre-trained LM may be sufficient, but can improve if topped with extra-training on in-domain data.

⁹Tokenized with the Twitter-specific “Twikenizer”: github.com/Guilherme-Routar/Twikenizer

Task	Tokens	RoB-Bs	RoB-Tw	% Diff
Emoji	14.3 \pm 7.4	22.4 \pm 7.4	21.6 \pm 6.8	2.8 \pm 6.9
Emotion	19.2 \pm 10.2	27.2 \pm 10.2	25.7 \pm 9.6	5.1 \pm 8.1
Hate	25.6 \pm 19.7	38.6 \pm 19.7	36 \pm 18.9	6.8 \pm 8.2
Irony	17.9 \pm 9.3	26.1 \pm 9.3	25.1 \pm 8.9	3.8 \pm 7.1
Sentiment	18.9 \pm 9.2	26.7 \pm 9.2	26.2 \pm 9.1	1.4 \pm 8.5
Offensive	28.4 \pm 20.9	41.9 \pm 20.9	39.4 \pm 19.7	5.7 \pm 8.5
Stance	20.6 \pm 7.1	30.7 \pm 7.1	30.5 \pm 6.9	0.5 \pm 4.8

Table 4: Tokenization statistics for all TWEETEVAL tasks. “Tokens” is the average number of tokens in each tweet using Twikenizer. RoB-RT and Rob-Tw refers to the average number of word pieces after tokenization with the original Roberta-base and with the model trained from scratch. “Diff” is the relative difference (%) of tokens in each tweet between these two tokenizers (if the difference is positive, the original RoBERTa includes more tokens). For stance detection, we computed the average statistics among the five targets.

For this initial benchmark and in the interest of reproducibility and accessibility, we focused on a fixed setting (i.e. classification). However, we acknowledge that other important tasks may need to be evaluated differently. Thus, for future work we would like to include more tasks in the context of social media NLP research. Potential improvements include, for example, accounting for the original multi-label nature of emotion classification, or covering more than only 20 emoji in emoji prediction. There are also other scenarios to be addressed as well, like sequence tagging (Baldwin et al., 2015; Gimpel et al., 2018), multimodality (Schifanella et al., 2016; Lu et al., 2018), and code-switching tasks (Barman et al., 2014; Vilares et al., 2016). This is similar to the evolution of GLUE (Wang et al., 2019b) into SuperGLUE (Wang et al., 2019a), with both benchmarks contributing to the development of the field in different ways. It is also important to highlight that these datasets do not represent their underlying tasks as a whole but only a subsample, and therefore contain biases - automatic models trained on them might not be able to generalize to other specific settings (Augenstein et al., 2017; Wiegand et al., 2019).

Finally, this benchmark could foster research in multitask learning. The fact that several similar tasks co-exist (e.g. sentiment analysis and emotion detection, or hate speech detection and offensive language identification) can lead to interesting analyses where the similarity of these tasks is exploited.

References

- Isabelle Augenstein, Leon Derczynski, and Kalina Bontcheva. 2017. Generalisation in named entity recognition: A quantitative analysis. *Computer Speech & Language*, 44:61–83.
- Timothy Baldwin, Paul Cook, Marco Lui, Andrew MacKinlay, and Li Wang. 2013. [How noisy social media text, how different social media sources?](#) In *Proceedings of the Sixth International Joint Conference on Natural Language Processing*, pages 356–364, Nagoya, Japan. Asian Federation of Natural Language Processing.
- Timothy Baldwin, Marie-Catherine de Marneffe, Bo Han, Young-Bum Kim, Alan Ritter, and Wei Xu. 2015. Shared tasks of the 2015 workshop on noisy user-generated text: Twitter lexical normalization and named entity recognition. In *Proceedings of the Workshop on Noisy User-generated Text*, pages 126–135.
- Francesco Barbieri, Jose Camacho-Collados, Francesco Ronzano, Luis Espinosa Anke, Miguel Ballesteros, Valerio Basile, Viviana Patti, and Horacio Saggion. 2018. Semeval 2018 task 2: Multilingual emoji prediction. In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 24–33.
- Utsab Barman, Amitava Das, Joachim Wagner, and Jennifer Foster. 2014. Code mixing: A challenge for language identification in the language of social media. In *Proceedings of the first workshop on computational approaches to code switching*, pages 13–23.
- Valerio Basile, Cristina Bosco, Elisabetta Fersini, Debora Nozza, Viviana Patti, Francisco Manuel Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. 2019. [SemEval-2019 task 5: Multilingual detection of hate speech against immigrants and women in twitter](#). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 54–63, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. [Enriching word vectors with subword information](#). *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Çağrı Çöltekin and Taraka Rama. 2018. Tübingen-oslo at semeval-2018 task 2: Svms perform better than rnns in emoji prediction. In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 34–38.
- Alexis Conneau and Douwe Kiela. 2018. [SentEval: An evaluation toolkit for universal sentence representations](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC-2018)*, Miyazaki, Japan. European Languages Resources Association (ELRA).
- Leon Derczynski, Alan Ritter, Sam Clark, and Kalina Bontcheva. 2013. [Twitter part-of-speech tagging for all: Overcoming sparse and noisy data](#). In *Proceedings of the International Conference Recent Advances in Natural Language Processing RANLP 2013*, pages 198–206, Hissar, Bulgaria. INCOMA Ltd. Shoumen, BULGARIA.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Kevin Gimpel, Nathan Schneider, Brendan O’Connor, Dipanjan Das, Daniel Mills, Jacob Eisenstein, Michael Heilman, Dani Yogatama, Jeffrey Flanigan, and Noah A. Smith. 2018. [Part-of-speech tagging for twitter: Annotation, features, and experiments](#).
- Bo Han and Timothy Baldwin. 2011. Lexical normalisation of short text messages: Makn sens a# twitter. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 368–378. Association for Computational Linguistics.
- Yuheng Hu, Kartik Talamadupula, and Subbarao Kambhampati. 2013. Dude, srsly?: The surprisingly formal nature of twitter’s language. In *Seventh International AAAI Conference on Weblogs and Social Media*.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. [Bag of tricks for efficient text classification](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain. Association for Computational Linguistics.
- Suin Kim, Ingmar Weber, Li Wei, and Alice Oh. 2014. Sociolinguistic analysis of twitter in multilingual societies. In *Proceedings of the 25th ACM conference on Hypertext and social media*, pages 243–248.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Di Lu, Leonardo Neves, Vitor Carvalho, Ning Zhang, and Heng Ji. 2018. Visual attention model for name tagging in multimodal social media. In *Proceedings of the 56th Annual Meeting of the Association for*

- Computational Linguistics (Volume 1: Long Papers)*, pages 1990–1999.
- Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. Semeval-2018 task 1: Affect in tweets. In *Proceedings of the 12th international workshop on semantic evaluation*, pages 1–17.
- Saif Mohammad, Svetlana Kiritchenko, Parinaz Sobhani, Xiaodan Zhu, and Colin Cherry. 2016. Semeval-2016 task 6: Detecting stance in tweets. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 31–41.
- Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, and Rada Mihalcea. 2020. Beneath the tip of the iceberg: Current challenges and new directions in sentiment analysis research. *arXiv preprint arXiv:2005.00357*.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding with unsupervised learning. *Technical report, OpenAI*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9.
- Alan Ritter, Sam Clark, Oren Etzioni, et al. 2011. Named entity recognition in tweets: an experimental study. In *Proceedings of the conference on empirical methods in natural language processing*, pages 1524–1534. Association for Computational Linguistics.
- Sara Rosenthal, Noura Farra, and Preslav Nakov. 2019. Semeval-2017 task 4: Sentiment analysis in twitter. *arXiv preprint arXiv:1912.00741*.
- Rossano Schifanella, Paloma de Juan, Joel Tetreault, and Liangliang Cao. 2016. [Detecting Sarcasm in Multimodal Social Platforms](#). *arXiv e-prints*, page arXiv:1608.02289.
- Cynthia Van Hee, Els Lefever, and Véronique Hoste. 2018. Semeval-2018 task 3: Irony detection in english tweets. In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 39–50.
- David Vilares, Miguel A Alonso, and Carlos Gómez-Rodríguez. 2016. En-es-es: An english-spanish code-switching twitter corpus for multilingual sentiment analysis. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 4149–4153.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019a. SuperGLUE: A stickier benchmark for general-purpose language understanding systems. In *Proceedings of NeurIPS*.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019b. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of ICLR*.
- Michael Wiegand, Josef Ruppenhofer, and Thomas Kleinbauer. 2019. [Detection of Abusive Language: the Problem of Biased Datasets](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 602–608, Minneapolis, Minnesota. Association for Computational Linguistics.
- Wei Xu. 2017. From shakespeare to twitter: What are language styles all about? In *Proceedings of the Workshop on Stylistic Variation*, pages 1–9.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages 5754–5764.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019. [SemEval-2019 task 6: Identifying and categorizing offensive language in social media \(OffenseEval\)](#). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 75–86, Minneapolis, Minnesota, USA. Association for Computational Linguistics.